

(12) NACH DEM VERTRAG ÜBER DIE INTERNATIONALE ZUSAMMENARBEIT AUF DEM GEBIET DES
PATENTWESENS (PCT) VERÖFFENTLICHTE INTERNATIONALE ANMELDUNG(19) Weltorganisation für geistiges Eigentum
Internationales Büro(43) Internationales Veröffentlichungsdatum
26. Februar 2004 (26.02.2004)

PCT

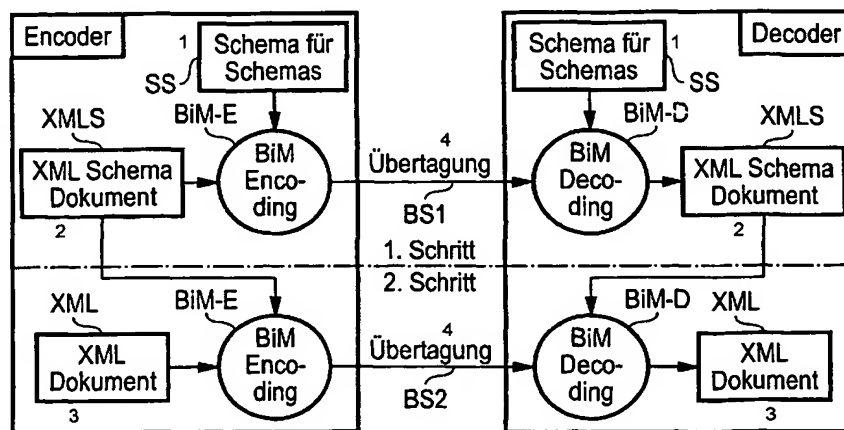
(10) Internationale Veröffentlichungsnummer
WO 2004/017225 A1

- (51) Internationale Patentklassifikation⁷: G06F 17/21, (72) Erfinder; und
17/22 (75) Erfinder/Anmelder (nur für US): HEUER, Jörg
[DE/DE]; Fischbachauerstr. 8, 81539 München (DE).
(21) Internationales Aktenzeichen: PCT/DE2003/002274 ✓ HUTTER, Andreas [DE/DE]; Sturmweg 42, 81673
München (DE); NIEDERMEIER, Ulrich [DE/DE];
(22) Internationales Anmeldedatum: 7. Juli 2003 (07.07.2003) ✓ Viehauserstrasse 18, 94405 Landau (DE).
- (25) Einreichungssprache: Deutsch (74) Gemeinsamer Vertreter: SIEMENS AKTIENGE-
SELLSCHAFT; Postfach 22 16 34, 80506 München (DE).
- (26) Veröffentlichungssprache: Deutsch
- (30) Angaben zur Priorität: 15. Juli 2002 (15.07.2002) DE (15 Jan 05)
102 31 971.5 18. Oktober 2002 (18.10.2002) DE
102 48 758.8
- (71) Anmelder (für alle Bestimmungsstaaten mit Ausnahme von
US): SIEMENS AKTIENGESELLSCHAFT [DE/DE];
Wittelsbacherplatz 2, 80333 München (DE).
- (81) Bestimmungsstaaten (national): AE, AG, AL, AM, AT,
AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR,
CU, CZ, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH,
GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC,
LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW,
MX, MZ, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC,

[Fortsetzung auf der nächsten Seite]

(54) Title: METHOD AND DEVICES FOR ENCODING/DECODING STRUCTURED DOCUMENTS, ESPECIALLY XML DOCUMENTS

(54) Bezeichnung: VERFAHREN UND VORRICHTUNGEN ZUM KODIEREN/DEKODIEREN VON STRUKTURIERTEN DOKUMENTEN, INSBESONDERE VON XML-DOKUMENTEN



1. SCHEMA FOR SCHEMAS
2. XML SCHEMA DOCUMENT
3. XML DOCUMENT
4. TRANSMISSION

(57) Abstract: The invention essentially relates to an encoding method for producing a bit stream or part of a bit stream from a schema according to a metaschema. According to the invention at least one of the following optimisation processes is carried out: separation of anonymous types from element declarations and attribute declarations, and encoding as own type, the type definition thereof as top level element being instantiated in the schema definition; normalisation of syntax trees on the encoder side; replacement of the character strings of type names; and transmission of information for the inheritance tree. The decoding takes said optimisation processes into account and conversely produces a schema from the bit stream.

[Fortsetzung auf der nächsten Seite]



SD, SE, SG, SK, SL, SY, TJ, TM, TN, TR, TT, TZ, UA,
UG, US, UZ, VC, VN, YU, ZA, ZM, ZW.

Veröffentlicht:

— mit internationalem Recherchenbericht

- (84) **Bestimmungsstaaten (regional):** ARIPO-Patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), eurasisches Patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), europäisches Patent (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PT, RO, SE, SI, SK, TR), OAPI-Patent (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

Zur Erklärung der Zweibuchstaben-Codes und der anderen Abkürzungen wird auf die Erklärungen ("Guidance Notes on Codes and Abbreviations") am Anfang jeder regulären Ausgabe der PCT-Gazette verwiesen.

(57) **Zusammenfassung:** Die Erfindung besteht im wesentlichen darin, mit einem Encodierverfahren aus einem Schema in Abhängigkeit eines Metaschemas einen Bitstrom oder einen Teil eines Bitstromes zu erzeugen, wobei eine oder mehrere der folgenden Optimierungen durchgeführt werden: - Abspaltung von anonymous Types aus Elementdeklarationen und Attributdeklarationen, und Codierung als eigener Typ, dessen Typdefinition als Top-Level Element in der Schema Definition instantiiert ist, - Normalisierung der Syntax Trees auf Encoderseite, - Ersetzung der Zeichenketten von Typnamen - Übertragung von Informationen für den Vererbungsbaum. Die Dekodierung berücksichtigt diese Optimierungen und erzeugt umgekehrt aus dem Bitstrom ein Schema.

Beschreibung

VERFAHREN UND VORRICHTUNGEN ZUM KODIEREN/DEKODIEREN VON STRUKTURIERTEN
DOKUMENTEN, INSBESONDERE VON XML-DOKUMENTEN

5

Die Erfindung betrifft Verfahren bzw. Vorrichtungen zum Enco-
dieren von strukturierten Dokumenten, insbesondere XML-
Dokumenten, bei denen aus einem strukturierten Dokument in
Abhängigkeit eines Schemas ein Bitstrom erzeugt wird und ein
10 Verfahren bzw. eine Vorrichtung zum Decodieren, bei denen aus
einem Bitstrom in Abhängigkeit eines Schemas ein strukturier-
tes Dokument erzeugt wird.

Im Rahmen der Arbeit am MPEG-7 Standard wurde ein Verfahren
15 zur binären Codierung von XML Daten entwickelt, das im fol-
genden BiM-Verfahren genannt wird und beispielsweise aus der
Veröffentlichung ISO/IEC FDIS 15938-1:2001(E), "Information
Technology - Multimedia Content Description Interface - Part
1: Systems" bekannt ist. Dieses Verfahren verwendet XML Sche-
20 ma Definitionen, die beim Encoder und Decoder vorliegen, bei-
spielsweise das MPEG-7 Schema, um die Codes für die einzelnen
Datenelemente der XML Beschreibung zu generieren. Dieses Ver-
fahren setzt voraus, dass dem Encoder und dem Decoder zumin-
dest teilweise die selben Schemadefinitionen vorliegen. Dies
25 kann beispielsweise gewährleistet werden, indem ein standar-
disiertes XML Schema im Decoder fest eingebaut wird. Außerdem
besteht die Möglichkeit das Schema separat oder zusätzlich
zum eigentlichen Dokument dem Decoder zu übermitteln. Die
Übertragung des Schemas vom Encoder zum Decoder kann in tex-
30 tueller Form durchgeführt werden, wobei eine Standard Text-
kompression, wie z.B. ZIP, angewendet werden kann.

Die der Erfindung zu Grunde liegende Aufgabe besteht nun dar-
in, Verfahren bzw. Vorrichtungen derart anzugeben, dass die
35 Übertragung des Schemas besonders effizient erfolgt und dass
die übertragene Datenmenge und die Rechenleistung am Decoder,
die für die Erzeugung der Codetabellen aus dem Schema nötig

ist, reduziert wird. Außerdem soll die Konsistenz eines nicht vollständig übertragenen Schemas sichergestellt werden.

5 Diese Aufgabe wird hinsichtlich des Encodierverfahrens durch die Merkmale des Patentanspruchs 1, hinsichtlich des Decodierverfahrens durch die Merkmale des Patentanspruchs 7, hinsichtlich der Encodiervorrichtung durch die Merkmale des Patentanspruchs 14 und hinsichtlich der Decodiervorrichtung durch die Merkmale des Patentanspruchs 15 erfindungsgemäß gelöst.

Die weiteren Ansprüche betreffen vorteilhafte Ausgestaltungen der erfindungsgemäßen Verfahren bzw. Vorrichtungen.

15 Die Erfindung besteht im wesentlichen darin, mit einem Encodierverfahren aus einem Schema in Abhängigkeit eines Metaschemas einen Bitstrom oder einen Teil eines Bitstromes zu erzeugen, wobei eine oder mehrere der folgenden Optimierungen durchgeführt werden:

20 - Abspaltung von anonymous Types aus Elementdeklarationen und Attributdeklarationen, und Codierung als eigener Typ, dessen Typdefinition als Top-Level Element in der Schema Definition instantiiert ist,

- Normalisierung der Syntax Trees auf Encoderseite,

25 - Ersetzung der Zeichenketten von Typnamen

- Übertragung von Informationen für den Vererbungsbaum.

Die Dekodierung berücksichtigt diese Optimierungen und erzeugt umgekehrt aus dem Bitstrom ein Schema.

30 Die Erfindung wird nachfolgend anhand von in den Zeichnungen dargestellten Ausführungsbeispielen erläutert. Dabei zeigt

Figur 1 eine Prinzipdarstellung zu Erläuterung der erfindungsgemäßen Encodierung/Decodierung,

35

Figur 2 eine Darstellung zur Erläuterung der Details einer bevorzugten Ausführungsform der Erfindung,

Figur 3 eine Darstellung zur Erläuterung der Details einer weiteren bevorzugten Ausführungsform der Erfindung und

5

Figur 4 eine Prinzipdarstellung einer bevorzugten Ausführungsform eines erfindungsgemäßen Decoders.

Da XML Schemas ihrerseits XML Dokumente sind, denen eine standardisierte Syntaxdefinition zugrunde liegt, nämlich ein sogenanntes "Schema for Schemas" (W3C Spezifikation), das quasi ein Metaschema darstellt, kann ein Schema ebenfalls mit dem oben genannten BiM-Verfahren codiert und übertragen werden.

15

In Figur 1 ist eine Anordnung gezeigt, bei der, in einem ersten Schritt, mit einem BiM-Encodierverfahren BiM-E aus einem XML-Schema XMLS in Abhängigkeit eines Metaschemas SS einen Teil eines Bitstromes oder einen Bitstrom BS1 erzeugt wird und bei der, in einem zweiten Schritt, mit dem selben BiM-Encodierverfahren BiM-E aus einem XML-Dokument XML in Abhängigkeit des Schemas XMLS ein weiterer Teil des Bitstromes oder ein Bitstrom BS2 erzeugt wird sowie in umgekehrter Richtung mit einem BiM-Decodierverfahren BiM-D aus den beiden

20 Teile des Bitstromes oder aus den Bitströmen BS1 und BS 2 ein XML Schema und das XML-Dokument wiedergewonnen werden.

25

In einer ersten bevorzugten Ausgestaltung der Erfindung wird eine Abspaltung von sogenannten „anonymous Types“ aus der Element- bzw. Attributdeklaration vorgenommen.

30

Die Übertragung eines XML Dokuments erfolgt beim BiM-Verfahren "depth first", der Vorgang der Schema Kompilierung am Decoder verlangt aber einen Aufbau "breadth first", wobei die diese Ausdrücke bspw. auf der Internetseite http://www.generation5.org/simple_search.shtml näher erläutert sind. Bei Gruppen wie Sequence oder Choice kann dies

35

durch einen kleinen Zwischenspeicher auf Decoderseite aus-
geglichen werden, aber bei den "anonymous Types", die den Typ
eines einzelnen Elements oder Attributs definieren können
rechtfertigt der Aufwand eine Umstrukturierung auf Encoder-
5 seite: die anonymous Type Definitionen, im nachfolgenden Bei-
spiel mit AT0 bezeichnet, werden aus der Elementdeklaration
des Elements " CurriculumVitae" herausgelöst und erhalten ei-
nen Namen und/oder Code, der zur Referenzierung bei dem ent-
sprechenden Element verwendet wird.

10 Vorteilhafterweise wird hierdurch die Tiefe der Hierarchie
der übertragenen Typen reduziert, wodurch die Kompilierung
des Schemas auf der Decoderseite vereinfacht wird.

Beispiel:

15 Schema vor der Umstrukturierung

```
<complexType name="PersonDescriptor">
  <element name="CurriculumVitae">
20   <complexType>
     <element name="name" type="string"/>
     <element name="birthday" type="date"/>
     ...
   </complexType>
25 </element>
   <element name="profession" type="profTp" />
</complexType>
```

30 Schema nach der Umstrukturierung

```
<complexType name="PersonDescriptor">
  <element name="CurriculumVitae" type="AT0"/>
  <element name="profession" type="profTp" />
35 </complexType>
```

5

```
<complexType name="AT0">  
  <element name="name" type="string"/>  
  <element name="birthday" type="date"/>  
  ...  
</complexType>
```

10 In einer zweiten bevorzugten Ausgestaltung der Erfindung wird die Normalisierung der Syntax Trees, wie sie in BiM spezifiziert ist, auf der Encoderseite durchgeführt.

15 Im BiM-Verfahren werden sogenannte „Finite State Automats“ die zur Dekodierung des Bitstroms verwendet werden aus Syntax Trees erzeugt, welche die Struktur des XML Schemas abbilden. Um die Codiereffizienz zu steigern entsprechen diese Syntax Trees nicht 1:1 den textuellen XML Definitionen, sondern es werden Normalisierungen vorgenommen. Drei verschiedene Fälle können hierbei auftreten:

20 1. Vereinfachung einer Gruppe, die nur ein Element enthält: Die Gruppe wird aufgelöst, und das enthaltene Element wird auf der Ebene der aufgelösten Gruppe in das Content Modell einsortiert, wobei die Attribute minOccurs und maxOccurs des Elements durch das Produkt der entsprechenden Attribute der aufgelösten Gruppe und des Elements vor der Umgruppierung ersetzt wird.

2. Vereinfachung einer choice-Gruppe, die ein Element mit dem Attributwert minOccurs=0 enthält:
30 Das Attribut „minOccurs“ der choice Gruppe wird unabhängig vom vorhergehenden Wert auf 0 gesetzt, das Element, das einen Attributwert minOccurs=0 hatte, wird einen Attributwert minOccurs=1 zugewiesen.

35

3. Vereinfachung von verschachtelten choice-Gruppen:

Enthält eine choice Gruppe eine andere choice Gruppe, die die Attributwerte minOccurs=maxOccurs=1 enthält, so wird diese choice Gruppe aufgelöst, und der Inhalt direkt der darüber-

5 liegenden choice Gruppe eingegliedert.

Diese Vereinfachungen sollten bei der Übertragung des Schemas schon am Encoder vorgenommen werden, da die Syntax Tree Transformationen die Vergabe der normativen Codes beeinflusst, und die Kompilierung des Schemas auf Decoderseite vereinfacht wird, wenn das Content-Modell direkt übernommen werden kann.

10

Die Vorteile liegen hier darin, dass hierdurch ebenfalls der Decoder entlastet wird und das Content-Modell direkt wie es bei der Typdecodierung entsteht dem Schema-Compiler zugeführt werden kann.

15

In einer dritten bevorzugten Ausgestaltung der Erfindung wird, wie in Figur 2 gezeigt, eine Ersetzung der Zeichenketten von Typnamen durchgeführt.

20

Im Attribut "name" und "base" einer Typdefinition, sowie beim Attribut "type" einer Element- oder Attributdeklaration treten häufig im Schema die selben Typnamen auf, die als Zeichenkette mehrfach übertragen werden würden. Bei der Codierung von Typnamen ist es deshalb vorteilhaft anstatt des Namens nur eine Nummer zu codieren, und separat dazu eine Tabelle, welche die Nummern wieder zu den ursprünglichen Namen in Beziehung setzt. Als Nummer bietet sich die Typnummer an, die der unten noch näher erläuterte Vererbungsbaum des Ur-Typs allen complexTypes zuordnet.

25

30

Entsprechendes gilt auch für das Attribut „name“ von globalen Element-Deklarationen und deren Referenzen in „ref“-Attributen und für den Namen von Ersetzungsgruppen im Attri-

35

but "substitutionGroup". In diesen Fällen kann beispielsweise der Schemaverzweigungscode SBC der globalen Elemente verwendet werden.

- 5 Hiermit kann Datenvolumen eingespart werden, da eine wiederholte Referenzierung auf den selben Typnamen kompakter dargestellt werden kann und die Typzuordnungstabelle mit einem Standardkompressor besser komprimiert werden kann, da die Typnamen nicht über den Bitstrom verteilt auftreten, sondern
10 kompakt in einem zusammenhängenden Bereich im Bitstrom.

In einer vorteilhaften Ausführungsform wird eine Liste umfassend die Typ- oder Elementnamen oder Namen von Erstetzungsgruppen codiert. Statt Nummern den Namen explizit zuzuordnen,
15 wird in dieser Ausführungsform die Position eines Namen in der Liste als Nummer verwendet. Dies ist vorteilhaft, da in der Liste keine Nummern mehr codiert werden müssen und somit eine effizientere Übertragung gewährleistet ist.

- 20 In einer vierten bevorzugten Ausgestaltung der Erfindung erfolgt eine Übertragung von Informationen für den Vererbungsbaum.

Jede Typdefinition enthält im sogenannten Attribut "base",
25 falls es vorhanden ist, die Information von welchem Typ er vererbt worden ist. Wenn alle diese Informationen für ein Schema gesammelt werden, ergibt sich eine Baumstruktur, der sogenannte Vererbungsbaum. Der Vererbungsbaum wird beim BiM-Codierungsverfahren verwendet, um im Falle einer Typumandlung
30 (type-cast) den neuen Typ des Elements zu übermitteln. Dabei ist der Code der allen vom Basistyp vererbten Typen zugeordnet wird, also der sogenannte Type Code, sowie die Länge dieses Codes für eine korrekte Dekodierung entscheidend. Die Länge ergibt sich aus der Gesamtzahl aller Typen im Vererbungsbaum unter dem Basistyp. Wenn das Schema vollständig
35 übertragen wurde lassen sich sowohl die Codes als auch die Codelänge auf der Decoderseite eindeutig ermitteln. Wenn aber

das Schema auf der Decoderseite nicht vollständig ist, muss noch Zusatzinformation übertragen werden, um bereits übertragenen Typen Type Codes zuzuweisen.

- 5 Jeder übertragene Typ hat im Namensfeld die Nummer des Typecode bezogen auf den Urtyp. Damit lässt sich der Typecode der abgeleiteten Typen durch einfache Differenzbildung ermitteln. Es fehlt noch die Information über die Mächtigkeit des durch den übertragenen Typen definierten Unterbaums, und damit die
10 Länge der Typecodes der von diesem übertragenen Typen abgeleiteten Typen. Diese Länge lässt sich mit wenigen Bits in einem variablen Längencode übertragen.

- In Figur 3 ist beispielhaft ein Vererbungsbaum eines Schemas
15 mit dem Typ A, von dem weitere Typen abgeleitet sind, dargestellt. Dieser Typ bekommt bezüglich des Urtyps "anyType" beispielsweise den Typecode 134. Von Typ A sind die Typen AA, AB und AC abgeleitet, deren Typecodes bezüglich des Urtyps angegeben sind. Um den Typecode bezüglich des Basistyps A zu
20 ermitteln, genügt es vom Typecode des gewünschten Typs den Typecode des Basistyps und eins zu subtrahieren:

$$TC_{Type} = TC_{Type \text{ bzgl. Urtyp}} - TC_{Basistyp \text{ bzgl. Urtyp}} - 1$$

- 25 Die fehlende Information über die Länge des Typecodes lässt sich am besten in der Referenztabelle als zusätzliche Zahl integrieren.

- Um die Information in der Typzuordnungstabelle mit einem
30 Standardkompressor komprimieren zu können empfiehlt es sich, sie auf ganze Bytes ausgerichtet abzulegen (bytealigned). Die erste Zahl ist eine vluimsbf5 Zahl, die die Zahl der Zeilen in der Tabelle codiert, dann folgt eine vluimsbf5 Zahl, die die Nummer an Bits für den Typecode codiert, und eine weitere
35 vluimsbf5 Zahl, die den Typecode bzgl. des Urtyps selbst darstellt. Es folgen Füllbits oder Stuffing Bits um die Ausrichtung auf Bytegrenzen zu erreichen.

Format der Typzuordnungstabelle			
Vuimsbf5	Vuimsbf5	Bits	Zeichen- kette
Zahl der Zeilen			
Länge Type- code 1	Typecode 1	0-7 Füll- bits	Name Typ 1
Länge Type- code 2	Typecode 2	0-7 Füll- bits	Name Typ 2
...

5 Die Übertragung einer Typzuordnungstabelle ermöglicht es, die
in einem kodierten Dokument evtl. vorhandenen Typecodes kor-
rekt zu decodieren, auch wenn das zugrundeliegende Schema
nicht oder noch nicht vollständig übertragen und/oder deco-
diert wurde.

10

Entsprechend sind mit globalen Elementen der globale SBC und
bei Elementen, die zu einer Ersetzungsgruppe gehören, der Er-
setzungscode zu übermitteln, wobei vorab für alle globalen
Elemente einmal die globale SBC-Länge und mit dem Kopfelement
15 der Ersetzungsgruppe die Länge des jeweiligen Ersetzungscode
übermittelt werden.

Es ist jede Kombination der in den einzelnen Ausgestaltungen
dargestellten Merkmale bei der Encodierung möglich und kann
20 in entsprechender Weise auch bei der Decodierung Eingang fin-
den.

Das BiM-Verfahren erfordert es, dass das XML-Schema in ein
Format kompiliert wird, das die Bestimmung der Länge der Co-
25 deworte und die Auswahl der Datenelemente durch die Werte der

Codes gestattet. Dafür gibt es mehrere Möglichkeiten. Im MPEG-7 Standard (ISO/IEC 15938-1:2001 Part1: Systems bzw. ISO/IEC 15938-6:2001 Part6: Referenzsoftware) ist für die Decodierung der Nutzlast bzw. Payload ein Modell vorgeschlagen, das Endliche Zustandsautomaten (Finite State Automats) verwendet, und für die Decodierung eines Context Pfades Codetabellen, die aus dem Schema generiert werden.

In einer in Figur 4 dargestellten bevorzugten Ausgestaltung des erfindungsgemäßen Decoders wird der Decodiervorgang mit einem Bytecodemodell beschrieben, wobei die Schemastruktur in ein System aus vernetzten Zuständen übersetzt wird, die von einem Bytecodeinterpreter BCI abgearbeitet werden, wobei ein vom Encoder empfangener Bitstrom BS die Information über den auszuwählenden Folgezustand enthält. Im Unterschied zu dem Modell, das im MPEG-7 Standard vorgeschlagen wird, ist das Bytecodemodell so angelegt, dass sowohl ein Bitstrom, der eine Payload repräsentiert, als auch ein Bitstrom der einen Context Pfad darstellt decodiert werden kann. Es ist deshalb nicht erforderlich dieselbe Information, die im Schema enthalten ist zweimal für die verschiedenen Codiervverfahren am Decoder vorzuhalten. Der Interpreter BCI liest die Information aus dem Eingangsbitstrom, die ein XML Dokument oder ein XML Schema im BiM Format codiert. Diese Information erlaubt die Auswahl unter den Folgezuständen des aktuellen Zustandes, der im Bytecode abgelegt ist. Die Folgezustände sind innerhalb des Bytecodes als Pointer P fest angelegt. Je nach Konfiguration wird ein Pfad, eine Payload oder ein Bytecode ausgegeben.

30

Die Decodierung eines Schemas läßt sich mit den oben vorgeschlagenen Modifikationen ebenfalls effizient im Bytecodemodell realisieren. In diesem Fall wird keine Payload und kein Pfad ausgegeben, sondern direkt Bytecode erzeugt, der vom Bytecodeinterpreter für die Decodierung der entsprechenden Typen verwendet werden kann.

35

Der Bytecode setzt sich aus Strukturelementen bzw. den Zuständen zusammen. Die Zustände sind von verschiedenem Typ, der mit dem Headerbitfeld des Zustandes identifiziert wird. Die Zustände enthalten abhängig vom Typ verschiedene Informationsfelder, die vom Bytecodeinterpreter gelesen, und je nach Konfiguration (Payload/ Context Pfad) und aktuellem Zustand ausgewertet werden.

Für die Arten von Zuständen, welche die Schemainformation repräsentieren sind mehrere Varianten denkbar. Wesentlich ist, dass sich durch die Zustände des Bytecodemodells alle Syntaxelemente eines XML Schemas nachbilden lassen, und dass die gesamte Information, die zur effizienten Decodierung der beiden im MPEG-7 Standard definierten Algorithmen (Context Pfad/Payload) notwendig ist, in den Zuständen zur Verfügung gestellt wird.

Ein möglicher Aufbau des Bytecodes wird im folgenden kurz dargestellt.

Arten von Zuständen, Übersicht:

1. Kopfzustand eines complexTyps

Der Kopfzustand eines Typs bildet den Einsprungspunkt bei der Decodierung eines complexType. Er enthält den Namen des Typs (falls es sich nicht um einen anonymen Typ handelt) sowie Information zu Vererbung des Typs (Zeiger auf Basiszustand) sowie Polymorphismus.

Spezifisch für die Payloadcodierung ist ein Zeiger auf eine Liste der Attribute des Typs. Spezifisch für die Context Pfad Codierung sind Felder mit der Zahl der Kindelemente für die Context - und Operand Tree Branch Code Tabellen.

Das letzte Informationsfeld ist ein Zeiger auf den Folgezustand, d.h. der erste Zustand, der den Inhalt des complexTypes repräsentiert (beispielsweise ein Elementzustand oder ein Auswahlzustand).

Graphische Darstellung eines Kopfzustands:

Headerbitfeld
Pointer auf String mit Name
Pointer auf Kopfzust. Basistyp
Pointer auf Vererbungsbaum
Zahl der Kinder Context TBC
Zahl der Kinder Operand TBC
Pointer auf Folgezustand

2. Auswahlzustand

Ein Auswahlzustand bildet eine choice Gruppe des XML Schemas nach. Der Auswahlzustand enthält im wesentlichen eine Pointerliste mit möglichen Folgezuständen. Um den tatsächlich ausgewählten Zustand zu bestimmen muß bei der Decodierung einer Payload der Bitstrom gelesen werden. Vom Auswahlzustand gibt es zwei Varianten: einen Startzustand, der in die verschiedenen möglichen Folgezustände verzweigt, sowie einen Endzustand, der die Auswahl wieder zusammenfaßt.

3. Elementzustand

Der Elementzustand bildet eine Elementdeklaration in einem complexType eines Schemas nach. Er enthält einen Pointer auf eine Zeichenkette mit dem Namen des Elements, sowie einen Pointer auf den Kopfzustand des Typs. Ferner ist evtl. Information über die Länge des Position Codes (nur für Pfad- Decodierung) und für Substitution Groups vorhanden.

4. Attributzustand

Ein Attributzustand bildet eine Attributdeklaration eines Schemas nach. Enthalten sind ein Pointer auf den Namen des Attributs, sowie ein Pointer auf den Kopfzustand des simple-Type des Attributs.

5. Occurrencezustand

Ein Occurrencezustand bildet die minOccurs und maxOccurs Attribute nach, die bei einem XML Schema z.B. bei einem Element oder einer Gruppe (choice, sequence, ...) auftreten können.

- 5 Er enthält einen Zeiger auf den Folgezustand, falls eine weitere Instanz des Elements oder der Gruppe auftritt, sowie einen Zeiger auf den Folgezustand, falls die letzte Instanz der Gruppe codiert wurde. Da bei XML Schemas die Möglichkeit besteht, daß ein Element sich selbst enthält (in der complexType
- 10 Definition des Elements, oder in einer noch tieferen Verschachtelung tritt das Element selbst wieder auf) kann auch ein Occurrencezustand gleichzeitig mehr als einmal aktiv sein. Deshalb ist ein Zeiger auf einen Stapel innerhalb des Occurrencezustands erforderlich, die den aktuellen Zustand
- 15 jeder aktiven Instanz des Occurrencezustands sichert..

6. Endzustand eines Typs

- Der Endzustand eines Typs enthält eine Zeigerliste mit allen Attributen dieses Typs. Sie ist bei der Decodierung eines
- 20 Pfades erforderlich, da in den Tree Branch Code Tabellen alle Attribute am Ende der Tabelle einsortiert werden. Beim Erreichen eines Endzustands verzweigt der Bytecodeinterpreter hierarchisch in das Element, das diesen Typ aufgerufen hat. Die entsprechende Information über das aufrufende Element muß im
- 25 Arbeitsspeicher des Bytecodeinterpreters abgelegt sein.

7. Kopfzustand eines simpleTypes

- Dieser Zustand steuert die Decodierung von Inhalt, d.h. er enthält einen Pointer auf einen Codec, der spezifisch Daten
- 30 des betreffenden Typs aus dem Bitstrom lesen und decodieren kann. Der Typ des Codecs ist in einem Informationsfeld spezifiziert.

- Die wesentlichen Vorteile des Bytecodemodells im Vergleich
- 35 zum Stand der MPEG-7 Referenzsoftware sind:

1. Die Schemainformation wird für beide Codierverfahren (Context Pfad / Payload) nur einmal am Decoder repräsentiert. Der größte Teil der Information in den Bytecodezuständen sind für beide Verfahren relevant. Ein kleinerer Teil ist spezifisch für jeweils eines der beiden Verfahren. Deshalb ist die Darstellung der Schemainformation am Decoder sehr kompakt.
2. Das Bytecodemodell stellt ein wohldefiniertes Datenformat für Schemainformation zur Verfügung, das sich z.B. auch zum Vorkompilieren und Abspeichern eignet (anstatt dem XML-Schema als Text).
3. Die Ausführung des Bytecodes durch einen Standardprozessor kann sehr schnell erfolgen, da das Bytecodemodell den Decodiervorgang sehr gut vorbereitet. Alle Information ist direkt im Zustand über Zeiger verfügbar, und muß nicht (wie in ISO/IEC 15938-6, Part 6: Referenzsoftware) zum Teil erst in Listen gesucht werden.
- Ein entsprechender Encoder kann auf die selbe Art und Weise realisiert werden, wobei er in der Weise invers ist, als dass die Zustände von der textuellen Repräsentation des strukturierten Dokuments gesteuert werden und die Zustandsübergänge die binäre Repräsentation generieren.

Patentansprüche

1. Verfahren zum Encodieren von strukturierten Dokumenten,
5 insbesondere XML-Dokumenten,
bei dem, in einem ersten Schritt, die Struktur des Schemas (XMLS) normiert wird, wobei Gruppen mit Elementen und/oder Attributen vereinfacht werden,
bei dem mit einem Encodierverfahren (BiM-E) aus dem normier-
10 ten Schema in Abhängigkeit eines Metaschemas (SS) ein Teil eines Bitstromes oder ein Bitstrom (BS1) erzeugt wird.
2. Verfahren nach Anspruch 1,
bei dem, in einem weiteren Schritt, mit dem selben Encodier-
15 verfahren (BiM-E) aus einem Dokument (XML) in Abhängigkeit des Schema (XMLS) ein weiterer Teil des Bitstromes oder ein weiterer Bitstrom (BS2) erzeugt wird.
3. Verfahren nach Anspruch 1 oder 2,
20 bei dem Elementdeklarationen und/oder Attributdeklarationen der Schemadefinition eines strukturierten Dokuments derart umstrukturiert werden, dass anonyme Typdefinitionen (AT0) aus den Elementdeklarationen und/oder Attributdeklarationen her-
ausgelöst werden und einen Namen und/oder Code erhalten, der
25 zur Referenzierung bei dem entsprechenden Element verwendet wird.
4. Verfahren nach einem der Ansprüche 1 bis 3,
bei dem anstatt Typnamen und/oder Elementnamen und/oder Namen
30 von Ersetzungsgruppen nur Nummern sowie eine oder mehrere Tabellen mit einer Zuordnung zwischen Nummern und Typnamen und/oder Elementnamen und/oder Namen von Ersetzungsgruppen codiert werden.
- 35 5. Verfahren nach einem der Ansprüche 1 bis 4,
bei dem eine oder mehrere Listen umfassend die Typnamen und/oder Elementnamen und/oder Namen von Ersetzungsgruppen sowie die Positionen der Typnamen und/oder Elementnamen

und/oder Namen von Ersetzungsgruppen in der Liste anstatt Typnamen und/oder Elementnamen und/oder Namen von Ersetzungsgruppen codiert werden.

- 5 6. Verfahren nach einem der vorhergehenden Ansprüche,
bei dem Informationen für den Vererbungsbaum von Typen, glo-
balen Elementen und/oder Ersetzungsgruppen codiert werden,
wobei jeder Typ durch eine Information über seinen Typcode
bezogen auf den Urtyp und der Länge aller Typcodes, die sich
10 auf den beschriebenen Typen beziehen, beschrieben wird
und/oder jedes globale Element durch die Länge des SBC und
einen SBC und/oder jedes Element in einer Ersetzungsgruppe
durch die Länge der Ersetzungscodes und einen Ersetzungscode
beschrieben wird.
- 15 7. Verfahren zum Decodieren von strukturierten Dokumenten,
insbesondere XML-Dokumenten,
bei dem, mit einem Decodierverfahren (BiM-D) aus einem Teil
eines Bitstromes oder aus einem Bitstrom (BS1) in Abhängig-
20 keit eines Metaschemas (SS) ein Schema (XMLS) erzeugt wird,
bei dem im Bitstrom festgestellt wird, ob die Struktur des
Schemas bereits normiert wurde, wobei Gruppen mit Elementen
und/oder Attributen vereinfacht wurden, und für diesen Fall
keine Normierung durchgeführt wird und
- 25 8. Verfahren nach Anspruch 7,
bei dem, in einem zweiten Schritt, mit dem selben Decodier-
verfahren (BiM-D) aus einem weiteren Teil des Bitstromes oder
einem weiteren Bitstrom (BS2) in Abhängigkeit des Schema
30 (XMLS) ein Dokument (XML) erzeugt wird.
9. Verfahren nach Anspruch 7,
bei dem, während der Decodierung des Schemas (XMLS), mit dem
selben Decodierverfahren (BiM-D) aus einem weiteren Teil des
35 Bitstromes oder einen weiteren Bitstrom (BS2) in Abhängigkeit
des bereits decodierten Teils des Schemas (XMLS) ein Dokument
(XML) erzeugt wird.

10. Verfahren nach einem der Ansprüche 7 bis 9,
bei dem Elementdeklarationen und/oder Attributdeklarationen
eines strukturierten Dokuments derart umstrukturiert werden,
5 dass anonyme Typen (AT0), denen zur Übertragung ein Name
und/oder ein Code zugewiesen wurde, in die jeweilige Element-
deklaration oder Attributdeklaration eingefügt werden, von
der der jeweilige anonyme Typ referenziert wird.
- 10 11. Verfahren nach einem der Ansprüche 7 bis 10,
bei dem aus dem Bitstrom Typnamen und/oder Elementnamen
und/oder Namen von Ersetzungsgruppen über Nummern sowie einer
oder mehrer Tabellen mit einer Zuordnung zwischen Nummern und
Typnamen und/oder Elementnamen und/oder Namen von Ersetzungs-
15 gruppen decodiert werden.
12. Verfahren nach einem der Ansprüche 7 bis 11,
bei dem aus dem Bitstrom Typnamen und/oder Elementnamen
und/oder Namen von Ersetzungsgruppen über eine oder mehrere
20 Listen umfassend die Typnamen und/oder Elementnamen und/oder
Namen von Ersetzungsgruppen sowie die Positionen der Typnamen
und/oder Elementnamen und/oder Namen von Ersetzungsgruppen in
der Liste decodiert werden.
- 25 13. Verfahren nach einem der Ansprüche 7 bis 12,
bei dem zunächst aus dem Bitstrom Informationen für einen
Vererbungsbaum von Typen und/oder globalen Elementen und/oder
Ersetzungsgruppen decodiert werden, wobei jeder Typ durch ei-
ne Information über seinen Typcode bezogen auf den Urtyp und
30 der Länge aller Typcodes, die sich auf den beschriebenen Ty-
pen beziehen, beschrieben wird
und/oder jedes globale Element durch die Länge des SBC und
einen SBC und/oder jedes Element in einer Ersetzungsgruppe
durch die Länge der Ersetzungscodes und einen Ersetzungscode
35 beschrieben wird.

14. Vorrichtung zum Encodieren von strukturierten Dokumenten,
insbesondere XML-Dokumenten,
bei der eine Encodiereinheit vorhanden ist,
die, in einem ersten Schritt, die Struktur des Schemas (XMLS)
5 normieren, wobei Gruppen mit Elementen und/oder Attributen
vereinfacht werden,
die aus dem normierten Schema in Abhängigkeit eines Metasche-
mas (SS) einen Teil eines Bitstromes oder einen Bitstrom
(BS1) erzeugen.
- 10 15. Vorrichtung zum Decodieren von strukturierten Dokumenten,
insbesondere XML-Dokumenten,
bei der eine Decodiereinheit vorhanden ist,
die, aus einem Teil eines Bitstromes oder aus einem Bitstrom
15 (BS1) in Abhängigkeit eines Metaschemas (SS) ein Schema er-
zeugt,
bei der im Bitstrom festgestellt wird, ob die Struktur des
Schemas (XMLS) bereits normiert wurden, wobei Gruppen mit
Elementen und/oder Attributen vereinfacht wurden, und für
20 diesen Fall keine Normierung durchgeführt wird.
16. Vorrichtung nach Anspruch 14,
bei der die Encodiereinheit einen konfigurierbaren Bytecode-
interpreter aufweist, der Informationen in einem Bytecode in-
25 terpretiert und der, abhängig von der Konfigurierung, aus dem
strukturierten Dokument basierend auf einem Bytecode einen
Code erzeugt, der einen Pfad oder eine Nutzlast repräsen-
tiert.
- 30 17. Vorrichtung nach Anspruch 15,
bei der die Decodiereinheit einen konfigurierbaren Bytecode-
interpreter aufweist, der durch Informationen aus dem Bit-
strom konfigurierbar ist und der, abhängig von der Konfigu-
rierung, aus dem Bitstrom basierend auf einem Bytecode einen
35 Pfad, eine Nutzlast oder einen Bytecode erzeugt.

FIG 1

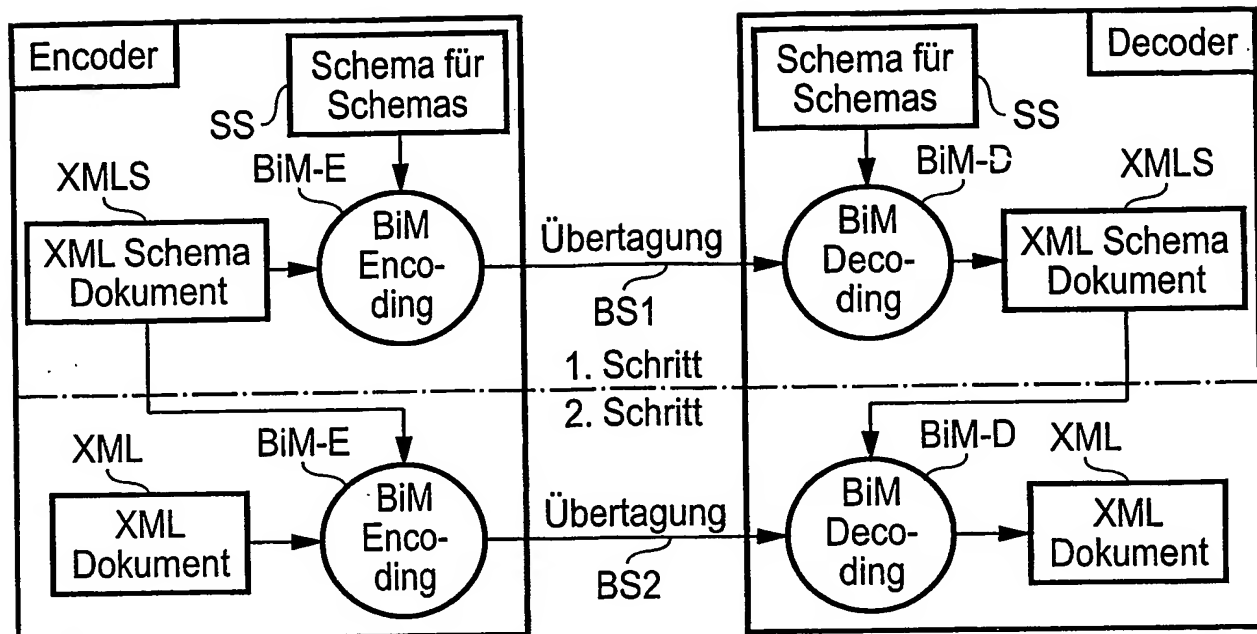


FIG 2

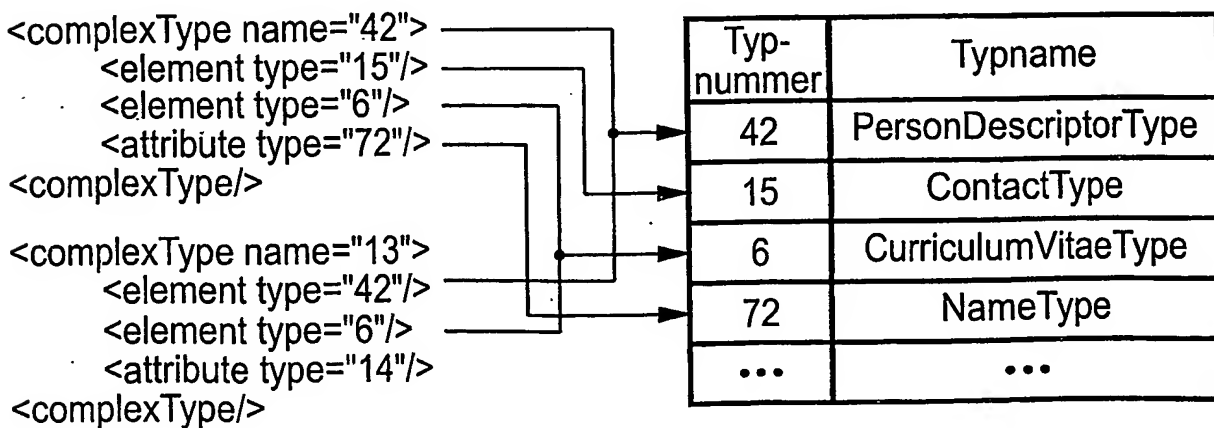


FIG 3

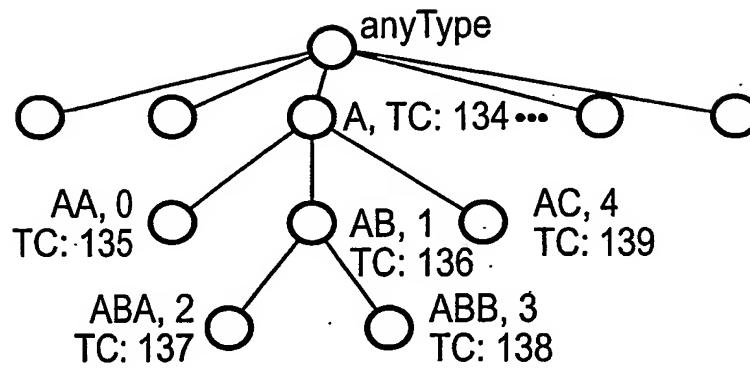
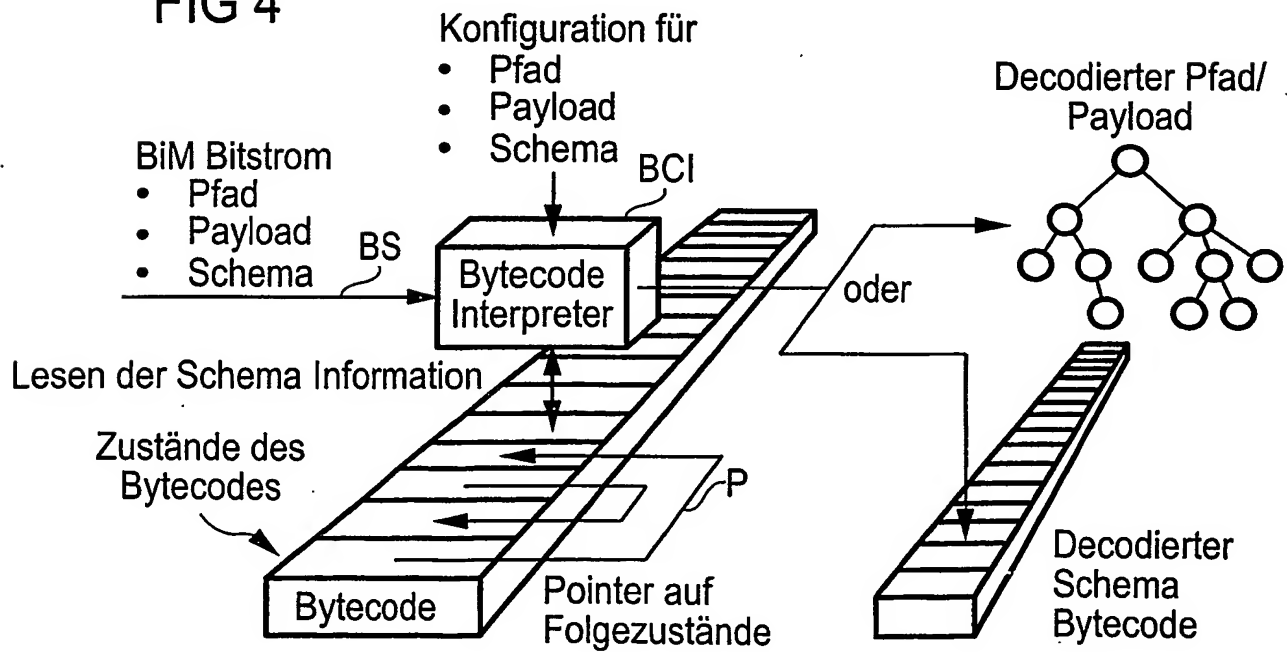


FIG 4



INTERNATIONAL SEARCH REPORT

International Application No

PCT/ 3/02274

A. CLASSIFICATION OF SUBJECT MATTER

IPC 7 G06F17/21 G06F17/22

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

IPC 7 G06F H04N H03M

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

EPO-Internal, WPI Data, PAJ, INSPEC

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	FR 2 813 743 A (SEYRAT CLAUDE) 8 March 2002 (2002-03-08) abstract page 1, line 1 - line 36 page 4, line 5 -page 6, line 7 page 9, line 15 -page 10, line 20; figures 1,2A-C,4 page 13, line 13 -page 14, line 9	1-3,5, 7-10,12, 14-17
Y	--- -/--	4,11,13

☒ Further documents are listed in the continuation of box C.☒ Patent family members are listed in annex.

* Special categories of cited documents:

- *A* document defining the general state of the art which is not considered to be of particular relevance
- *E* earlier document but published on or after the international filing date
- *L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- *O* document referring to an oral disclosure, use, exhibition or other means
- *P* document published prior to the international filing date but later than the priority date claimed

- *T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
- *X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
- *Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.
- *&* document member of the same patent family

Date of the actual completion of the international search

21 November 2003

Date of mailing of the international search report

02/12/2003

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Woods, J

INTERNATIONAL SEARCH REPORT

International Publication No
PCT 03/02274

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category °	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	"TEXT OF ISO/IEC FCD 15938-1 INFORMATION TECHNOLOGY - MULTIMEDIA CONTENT DESCRIPTION INTERFACE - PART 1 SYSTEMS" ISO/IEC JTC1/SC29/WG11 MPEG01/N4001, XX, XX, March 2001 (2001-03), pages 1-2, I-V, 6-58, XP001001465 cited in the application page 26, line 11 -page 28, line 9 page 34, line 26 -page 41, line 12; figures 11-13	1,3,5,6
Y	-----	13
Y	GIRARDOT M ET AL: "MILLAU: AN ENCODING FORMAT FOR EFFICIENT REPRESENTATION AND EXCHANGE OF XML OVER THE WEB" COMPUTER NETWORKS AND ISDN SYSTEMS, NORTH HOLLAND PUBLISHING. AMSTERDAM, NL, vol. 33, no. 1-6, June 2000 (2000-06), pages 747-765, XP001005949 ISSN: 0169-7552 page 750, left-hand column, line 5 -page 751, right-hand column, line 22; tables 1,2	4,11

INTERNATIONAL SEARCH REPORT

Information on patent family members

International Application No

PCT/03/02274

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
FR 2813743 A	08-03-2002	FR 2813743 A1	08-03-2002
		AU 8779601 A	22-03-2002
		EP 1316220 A1	04-06-2003
		WO 0221848 A1	14-03-2002

INTERNATIONALER RECHERCHENBERICHT

Internationales Patentzeichen

PCT/03/02274

A. KLASSIFIZIERUNG DES ANMELDUNGSGEGENSTANDES
IPK 7 G06F17/21 G06F17/22

Nach der Internationalen Patentklassifikation (IPK) oder nach der nationalen Klassifikation und der IPK

B. RECHERCHIERTE GEBIETE

Recherchierter Mindestprüfstoff (Klassifikationssystem und Klassifikationssymbole)
IPK 7 G06F H04N H03M

Recherchierte aber nicht zum Mindestprüfstoff gehörende Veröffentlichungen, soweit diese unter die recherchierten Gebiete fallen

Während der internationalen Recherche konsultierte elektronische Datenbank (Name der Datenbank und evtl. verwendete Suchbegriffe)

EPO-Internal, WPI Data, PAJ, INSPEC

C. ALS WESENTLICH ANGESEHENE UNTERLAGEN

Kategorie*	Bezeichnung der Veröffentlichung, soweit erforderlich unter Angabe der in Betracht kommenden Teile	Betr. Anspruch Nr.
X	FR 2 813 743 A (SEYRAT CLAUDE) 8. März 2002 (2002-03-08) Zusammenfassung Seite 1, Zeile 1 - Zeile 36 Seite 4, Zeile 5 -Seite 6, Zeile 7 Seite 9, Zeile 15 -Seite 10, Zeile 20; Abbildungen 1,2A-C,4 Seite 13, Zeile 13 -Seite 14, Zeile 9	1-3,5, 7-10,12, 14-17
Y	--- -/--	4,11,13

☒ Weitere Veröffentlichungen sind der Fortsetzung von Feld C zu entnehmen

☒ Siehe Anhang Patentfamilie

* Besondere Kategorien von angegebenen Veröffentlichungen :

A Veröffentlichung, die den allgemeinen Stand der Technik definiert, aber nicht als besonders bedeutsam anzusehen ist

E älteres Dokument, das jedoch erst am oder nach dem internationalen Anmeldedatum veröffentlicht worden ist

L Veröffentlichung, die geeignet ist, einen Prioritätsanspruch zweifelhaft erscheinen zu lassen, oder durch die das Veröffentlichungsdatum einer anderen im Recherchenbericht genannten Veröffentlichung belegt werden soll oder die aus einem anderen besonderen Grund angegeben ist (wie ausgeführt)

O Veröffentlichung, die sich auf eine mündliche Offenbarung, eine Benutzung, eine Ausstellung oder andere Maßnahmen bezieht

P Veröffentlichung, die vor dem internationalen Anmeldedatum, aber nach dem beanspruchten Prioritätsdatum veröffentlicht worden ist

T Spätere Veröffentlichung, die nach dem internationalen Anmeldedatum oder dem Prioritätsdatum veröffentlicht worden ist und mit der Anmeldung nicht kollidiert, sondern nur zum Verständnis des der Erfindung zugrundeliegenden Prinzips oder der ihr zugrundeliegenden Theorie angegeben ist

X Veröffentlichung von besonderer Bedeutung; die beanspruchte Erfindung kann allein aufgrund dieser Veröffentlichung nicht als neu oder auf erfinderischer Tätigkeit beruhend betrachtet werden

Y Veröffentlichung von besonderer Bedeutung; die beanspruchte Erfindung kann nicht als auf erfinderischer Tätigkeit beruhend betrachtet werden, wenn die Veröffentlichung mit einer oder mehreren anderen Veröffentlichungen dieser Kategorie in Verbindung gebracht wird und diese Verbindung für einen Fachmann naheliegend ist

& Veröffentlichung, die Mitglied derselben Patentfamilie ist

Datum des Abschlusses der internationalen Recherche

21. November 2003

Absenddatum des internationalen Recherchenberichts

02/12/2003

Name und Postanschrift der Internationalen Recherchenbehörde
Europäisches Patentamt, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Bevollmächtigter Bediensteter

Woods, J

C.(Fortsetzung) ALS WESENTLICH ANGESEHENE UNTERLAGEN

Kategorie°	Bezeichnung der Veröffentlichung, soweit erforderlich unter Angabe der in Betracht kommenden Teile	Betr. Anspruch Nr.
X	"TEXT OF ISO/IEC FCD 15938-1 INFORMATION TECHNOLOGY - MULTIMEDIA CONTENT DESCRIPTION INTERFACE - PART 1 SYSTEMS" ISO/IEC JTC1/SC29/WG11 MPEG01/N4001, XX, XX, März 2001 (2001-03), Seiten 1-2,I-V,6-58, XP001001465 in der Anmeldung erwähnt Seite 26, Zeile 11 -Seite 28, Zeile 9 Seite 34, Zeile 26 -Seite 41, Zeile 12; Abbildungen 11-13	1,3,5,6
Y	-----	13
Y	GIRARDOT M ET AL: "MILLAU: AN ENCODING FORMAT FOR EFFICIENT REPRESENTATION AND EXCHANGE OF XML OVER THE WEB" COMPUTER NETWORKS AND ISDN SYSTEMS, NORTH HOLLAND PUBLISHING. AMSTERDAM, NL, Bd. 33, Nr. 1-6, Juni 2000 (2000-06), Seiten 747-765, XP001005949 ISSN: 0169-7552 Seite 750, linke Spalte, Zeile 5 -Seite 751, rechte Spalte, Zeile 22; Tabellen 1,2 -----	4,11

INTERNATIONALER RECHERCHENBERICHT

Angaben zu Veröffentlichungen, die derselben Patentfamilie gehören

Internationale Zeichen

PCT/03/02274

Im Recherchenbericht angeführtes Patentdokument		Datum der Veröffentlichung	Mitglied(er) der Patentfamilie		Datum der Veröffentlichung
FR 2813743	A	08-03-2002	FR	2813743 A1	08-03-2002
			AU	8779601 A	22-03-2002
			EP	1316220 A1	04-06-2003
			WO	0221848 A1	14-03-2002
<hr/>					